

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 06-175787

(43)Date of publication of application : 24.06.1994

(51)Int.Cl.

G06F 3/06

G06F 12/00

(21)Application number : 04-330739

(71)Applicant : HITACHI LTD
HITACHI COMPUTER
PERIPHERALS CO LTD

(22)Date of filing : 10.12.1992

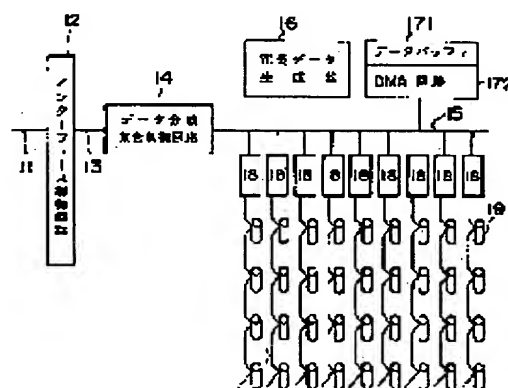
(72)Inventor : IWASAKI HIDEHIKO
TOKIDA YOSHINORI
MIZUNO KATSUTOSHI
SUZUKI RYOICHI
BABA HIDEMI

(54) DISK ARRAY DEVICE

(57)Abstract:

PURPOSE: To shorten the transfer time of redundant data to a disk device by generating redundant data in a unit smaller than a data distribution unit.

CONSTITUTION: Data inputted to a data distributing cluster control circuit 14 is divided into 512 bytes as a minimum distribution unit and is distributed in the 2n-fold distribution unit. A redundant data generator 16 manages distributed data in 512-byte unit by a higher-order data management counter, and generated data is managed in 512-byte units also by a redundant data management counter, and an input/output control circuit controls permission/inhibition of data input to the redundant data generator 16 while referring to both of them. An input/output permission/inhibition signal as the result is transmitted to a DMA circuit 172 in each 512-byte unit. The DMA circuit 172 refers to transmitted information to determine the service of data transfer from a disk control circuit.



LEGAL STATUS

[Date of request for examination] 13.01.1999

[Date of sending the examiner's decision of rejection] 26.09.2000

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3186272

[Date of registration] 11.05.2001

[Number of appeal against examiner's decision
of rejection]

[Date of requesting appeal against examiner's
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19)日本国特許庁(JP)

(12) 公開特許公報(A)

(11)特許出願公開番号

特開平6-175787

(43)公開日 平成6年(1994)6月24日

(51)Int.Cl. ⁵	識別記号	庁内整理番号	FI	技術表示箇所
G 0 6 F 3/06	3 0 1 Z	7165-5B		
12/00	5 4 5 A	8526-5B		

審査請求 未請求 請求項の数8(全 11 頁)

(21)出願番号 特願平4-330739

(22)出願日 平成4年(1992)12月10日

(71)出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(71)出願人 000233033

日立コンピュータ機器株式会社

神奈川県小田原市国府津2880番地

(72)発明者 岩崎 秀彦

神奈川県小田原市国府津2880番地 株式会

社日立製作所ストレージシステム事業部内

(72)発明者 常田 義則

神奈川県小田原市国府津2880番地 日立コ

ンピュータ機器株式会社内

(74)代理人 弁理士 富田 和子

最終頁に続く

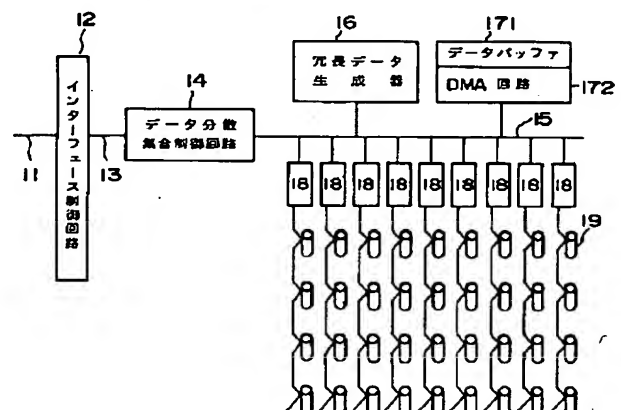
(54)【発明の名称】 ディスクアレイ装置

(57)【要約】

【目的】 ディスクアレイ装置において、冗長データ生成器へのデータ入力、及び、冗長データ生成器からのデータ出力の転送効率を向上する。

【構成】 冗長データ生成器16は、データ分散、集合制御回路14、及び、ディスク制御回路18からのデータを、分散単位よりも小さい単位にて入出力の可／不可制御をそれぞれ独立に行い、その制御結果をDMA回路172に伝達する。DMA回路172は、これを元に、データ分散集合制御回路14、及び、ディスク制御回路18に対するデータ転送のサービスを行う。

ディスクアレイ装置ブロック図(図1)



【特許請求の範囲】

【請求項1】複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、

上記冗長データ生成手段は、冗長データを生成する際、データ分散単位よりも小さい単位で、冗長データを生成することを特徴とするディスクアレイ装置。

【請求項2】複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、

上位装置、及びディスク装置ユニットの両方から冗長データ生成手段にデータを転送して、冗長データを生成する場合に、上位装置から転送されてくるデータ量をカウントする上位側データ管理手段と、

ディスク装置ユニットから転送されてくるデータ量をカウントするディスク側データ管理手段と、

冗長データ生成手段が生成した冗長データ量をカウントする冗長データ管理手段と、

上記上位側データ管理手段と上記ディスク側データ管理手段と上記冗長データ管理手段とがカウントした結果に従って、上記冗長データ生成手段へのデータ転送を制御する入出力制御手段とを有することを特徴とするディスクアレイ装置。

【請求項3】複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、

データを上位装置に転送する時に、冗長データ及び冗長データを生成するために使用した上記データにより、誤りの有無を調べる誤り検出手段を有することを特徴とするディスクアレイ装置。

【請求項4】複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、

データをディスク装置ユニットに書き込む際に、書き込むデータの量が分散単位よりも小さいときに、

上記ディスク装置ユニットは、分散単位中の書替の対象となる記憶位置にあるデータはそのまま転送し、分散単位中の書替の対象とならない記憶位置にあるデータは転送せず、

上記冗長データ生成手段は、上記ディスク装置ユニットから転送されてきたデータ及び書き込むデータから冗長データを生成することを特徴とするディスクアレイ装

置。

【請求項5】複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、

上記データを記録する時に、上記データに記録される位置を示す分散パラメータを付加する分散パラメータ付加手段を有し、

上記ディスク装置ユニットは、上記データとともに分散パラメータを記録することを特徴とするディスクアレイ装置。

【請求項6】請求項5記載のディスクアレイ装置において、

上記データを再生する時に、データと共に分散パラメータを各々のディスク装置ユニットより読みだし、分散パラメータが読みだした位置と整合が取れているかどうかを確認する分散パラメータ確認手段を有することを特徴とするディスクアレイ装置。

【請求項7】請求項5または6記載のディスクアレイ装置において、

上記冗長データ生成手段は、上記データより冗長データを生成する際、分散パラメータに対しても冗長データを生成し、分散パラメータを記録したディスク装置ユニットに障害が発生した時は、正常なディスク装置ユニットのデータと上記冗長データより、分散パラメータを再生成することを特徴とするディスクアレイ装置。

【請求項8】請求項5、6または7記載のディスクアレイ装置において、

上記データの記録される位置を示す位置情報を保持する位置保持手段と、

上記位置保持手段が保持する位置情報が消失した場合、ディスク装置ユニットに記録された分散パラメータにより位置情報を生成する位置情報生成手段とを有することを特徴としたディスクアレイ装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、複数台のディスク装置ユニットにデータを分散し、その分散されたデータに対して、冗長データを生成するディスクアレイ装置に関する。

【0002】

【従来の技術】ディスクアレイ装置に関する従来技術としては、例えば、特開平2-236714号公報「アレイ型ディスク駆動機構システム及び方法」等で示されるように、データを複数台のディスク装置ユニットに分散して記憶し、高速転送を実現するものや、分散データより冗長データを生成し、ディスク障害時、冗長データよりデータを再構築するものがある。

【0003】

【発明が解決しようとする課題】冗長データを有するディスクアレイ装置では、データライト時の冗長データの生成、及び生成した冗長データのディスク装置への転送が必要となり、このことがデータ転送効率を低下させる一因となる。

【0004】本発明の第1の目的は、冗長データのディスク装置への転送時間を短縮したディスクアレイ装置を提供することである。

【0005】また、データライト時に生成した冗長データは、データリード時には、リードされず、ディスク障害発生時のデータ再構築時にリードされる。このため、データリードを主に行われるディスクアレイ装置においては、冗長データ用ディスクへのアクセスが少ないため、障害検出が遅れ、2台のディスク装置ユニットが障害を起こす場合がある。2台のディスク装置ユニットが障害を起こすと原理的にデータ再構築不能となる。

【0006】本発明の第2の目的は、上記問題点を解決するため、障害検出の遅れを防ぎ、信頼性を向上させたディスクアレイ装置を提供することにある。

【0007】また、ディスクアレイ装置におけるデータは、冗長データを有することで、ディスク障害時等の信頼性を確保しているが、データの分散および集合はディスクアレイ装置内の回路が実行する為、分散、集合させる回路そのものに障害があった場合、または、分散記録した時のパラメータ（分散させたディスク装置ユニットの番号等）が、分散、集合させる回路内より消失した場合、データ障害の発生がなくともデータ破壊となる。

【0008】本発明の第3の目的は、分散、集合に対する信頼性を向上させたディスクアレイ装置を提供することにある。

【0009】

【課題を解決するための手段】上記第1の目的を達成するために、複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、上記冗長データ生成手段は、冗長データを生成する際、データ分散単位よりも小さい単位で、冗長データを生成することとした。

【0010】上記第2の目的を達成するために、複数のディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、データを上位装置に転送する時に、冗長データ及び冗長データを生成するために使用した上記データにより、誤りの有無を調べる誤り検出手段を有することとした。

【0011】上記第3の目的を達成するために、複数の

ディスク装置ユニットを有し、各々のディスク装置ユニットに上位装置から転送されてくるデータをデータ分散単位で分散して、記録するとともに、上記データに対して冗長データを生成する冗長データ生成手段を有するディスクアレイ装置において、上記データを記録する時に、上記データに記録される位置を示す分散パラメータを付加する分散パラメータ付加手段を有し、上記ディスク装置ユニットは、上記データとともに分散パラメータを記録することとした。

【0012】

【作用】ディスクアレイ装置において、データライト時、冗長データ生成手段により、冗長データを生成する際、データ分散単位よりも小さい単位で、冗長データを生成する。

【0013】冗長データ生成手段において、上位装置からのデータはデータ分散単位を分割した単位で入力され、また、生成した冗長データは、同様にデータ分散単位を分割した単位で出力制御される。このため、冗長データは、分散単位分の冗長データが生成されるのを待つことなく、1つの分割単位の冗長データが生成される毎に出力し、また、上位装置からのデータは、データ分散単位×分散ディスク装置数のデータ転送を行った後、生成した分散単位の冗長データが全て冗長データ生成手段より出力されるのを待つことなく、生成単位の冗長データが出力された時点で次の分散データの入力を行うことができる。

【0014】また、ディスクアレイ装置において、誤り検出手段は、データを上位装置に転送する時に、冗長データ及び冗長データを生成するために使用した上記データにより、誤りの有無を調べる。

【0015】このため、データリード時に冗長データも含めてリードし、データ照合するため、早い時点で冗長データが記録されたディスク装置ユニットの故障が検出でき、ディスク装置ユニットが2台故障することを防げる。従って、信頼性が向上する。

【0016】さらに、ディスクアレイ装置において、分散パラメータ付加手段は、上記データを記録する時に、上記データに記録される位置を示す分散パラメータを付加し、上記ディスク装置ユニットは、上記データとともに分散パラメータを記録する。

【0017】このため、データ記録時、分散パラメータが生成、付加され、その結果、分散先のディスク装置ユニット毎に、付加された分散パラメータを照合し、分散データの妥当性を保証して分散データと共にディスク装置に記録する等のことが可能となる。

【0018】

【実施例】以下、添付の図面に示す実施例により、更に詳細に本発明について説明する。図1は、本発明を説明するための冗長データを有するディスクアレイ装置の第1の実施例である。また、図2は、冗長データ生成手段

および誤り検出手段である冗長データ生成器の入出力制御ブロック図である。

【0019】本ディスクアレイ装置は、ホストSCSIバス11と、インターフェース制御回路12と、データバス13と、データ分散集合制御回路14と、共通データバス15と、データバッファ171と、DMA回路172と、冗長データ生成器16と、ディスク制御回路18と、ディスク装置ユニット（ディスク駆動機構）19とを有する。冗長データ生成器16は、上位側データ管理手段である上位側データ管理カウンタ21と、ディスク側データ管理手段であるディスク側データ管理カウンタ22と、冗長データ管理手段である冗長データ管理カウンタ23と、入出力制御手段である入出力制御回路24とを有する。

【0020】本ディスクアレイ装置は、コンピュータの記憶装置として動作し、汎用I/OインターフェースであるホストSCSIバス11を介して、コンピュータと接続される。

【0021】コンピュータのデータ書き込み要求時、ホストSCSIバス11を介して転送されるデータは、インターフェース制御回路12を経由し、データバス13よりデータ分散集合制御回路14に入力され、ここで分散される。分散されたデータは、DMA回路172により、共通データバス15を介して、ディスク制御回路18に解放されたデータバッファ171のエリアに転送される。この時、分散されたデータは、データバッファ171に転送すると同時に、冗長データ生成器16にも転送される。転送された各々の分散データは、排他ORが実行され、冗長データ生成器16内に冗長データが生成される。冗長データの生成が完了すると、分散データと同様にして該当するデータバッファ171のエリアに転送される。各々のディスク制御回路18は、自身に解放されたデータバッファ171のエリアより、指定されたディスク装置ユニット19に対して、分散されたデータ、及び、分散データより生成された冗長データを書込む。

【0022】コンピュータのデータ読みだし要求時は、該当するディスク制御回路18が、指定されたディスク装置ユニット19よりデータを読みだし、データバッファ171の定められたエリアに転送すると同時に、冗長データ生成器16にも転送する。冗長データ用のディスク制御回路18は、冗長データ生成器16にのみ転送する。データ分散集合制御回路14は、データバッファ171に転送された分散データを集合し、集合されたデータは、データバス13、インターフェース制御回路12を経由し、ホストSCSIバス11を介しコンピュータに転送される。また、冗長データ生成器16は、入力された分散データ、及び冗長データの排他ORを実行し、その結果がオール“00”であることを確認し、分散データと冗長データの整合性をチェックする。

【0023】分散データから生成され、ディスク装置ユニット19に書込まれた冗長データは、分散を行った複数のディスク装置ユニット19の中で、あるディスク装置ユニット19に障害が発生した場合の、障害が発生したディスク装置ユニット19の分散データを復元するために使用される。分散データの復元は、障害の発生したディスク装置ユニット19以外のディスク装置ユニット19より、分散されたデータ、及び冗長データを読みだし、冗長データ生成器16へ転送（分散データは、該当するデータバッファエリア19にも転送）することで行う。

【0024】冗長データ生成器16は、転送されたデータの排他ORを実行し、障害分散データを再構築する。再構築されたデータは、該当するデータバッファ171のエリアに転送され、データ分散集合制御回路14は、通常の読みだし処理と同様に、データバッファ171の分散データを集合し、集合されたデータは、コンピュータへと転送される。また、ディスクアレイ装置内にスベアディスク装置ユニットを備えている場合は、ディスクアレイ装置内で、障害ディスク装置ユニットの内容をスベアディスク装置ユニットに再構築する。

【0025】コンピュータのデータ書き込み要求時で、各々のディスク装置ユニット列の途中の列からデータ書き込みが開始され、なおかつ、その分散単位の途中からデータ書き込みが開始される場合について、更に詳細に説明する。この場合、コンピュータから転送されるデータからのみでは冗長データは生成できないため、ディスク装置ユニット19からのデータ読みだしも行わなければならない。

【0026】ホストSCSIバス11、インターフェース制御回路12、データバス13を介し、データ分散集合制御回路14に入力されたデータは、ここで最小分散単位である512バイトに分割され、なおかつ、その2n倍の分散単位にて分散される。

【0027】また、この時分散単位の途中からデータ書き込み開始となるため、データ分散集合制御回路14は、分散単位の開始点から書き込み開始までを“00”データを付加し転送する。分散されたデータは、DMA回路172によりデータバッファ172、及び、冗長データ生成器16に、512バイト単位で転送される。但し、付加された“00”データは、冗長データ生成器16にのみ転送される。この時、冗長データ生成器16は、上位側データ管理カウンタ21により、分散データを512バイトの単位に管理し、また、生成した冗長データについても、冗長データ管理カウンタ23により、512バイト単位に管理し、双方を参照しながら、入出力制御回路24は、冗長データ生成器16へのデータ入力の可／不可制御を行う。

【0028】その結果である入出力可／不可信号25は、DMA回路172に対し、512バイト単位毎に伝

達される。

【0029】DMA回路172は、その伝達された情報を参照し、データ分散集合制御回路14からのデータ転送のサービスを決定する。一方、該当するディスク制御回路18は、指定されたディスク装置ユニット19よりデータを読みだし、512バイト単位に分割する。

【0030】また、この時、分散単位の途中からデータ書き込み開始となるディスク制御回路18は、分散単位の開始点から書き込み開始点までを、“00”データを付加し転送する。分割されたデータは、DMA回路172により、冗長データ生成器16に転送されるこの時、冗長データ生成器16は、ディスク側データ管理カウンタ22により、ディスク制御回路18からのデータを、ディスク制御回路18毎に512バイトの単位に管理し、また、生成した冗長データについても、冗長データ管理カウンタ23により、512バイト単位に管理し、双方を参照しながら、冗長データ生成器16へのデータ入力の可／不可制御を行う。

【0031】また、上位側データ管理カウンタ21により、データ分散集合制御回路14からのデータも合わせて管理しており、冗長データ生成器16は、入力されるデータ全てを独立に制御するため、データ転送にソフト制御を必要としない。

【0032】冗長データ生成器16へのデータ入力の可／不可制御の結果は、入出力可／不可信号25として、DMA回路172に対し、512バイト単位毎に伝達される。DMA回路172は、その伝達された情報を参照し、ディスク制御回路18からのデータ転送のサービスを決定する。

【0033】データ分散集合制御回路14、及び、ディスク制御回路18から冗長データ生成器16に入力されたデータは、配置ORが実行され、冗長データ生成器16内に、新しい冗長データが生成される。この時、分散データが欠けている部分についても、“00”のダミーデータが転送されているため、排他ORを行う列数は、分散データが欠けている部分、揃っている部分共に同一列数として扱う。また、分散データが欠けている部分については、“00”データにて排他ORを実行するため、この部分については、排他ORを実行しないことと同様であり、ダミーデータを含んでも、冗長データは破壊されることはない。冗長データの生成が完了すると、該当するデータバッファ171のエリアに転送される。各々のディスク制御回路18は、自身に解放されたデータバッファ171のエリアより、指定されたディスク装置ユニット19に対して、分散されたデータ、及び、新しく生成された冗長データを書込む。

【0034】以上のように、本発明は冗長データを有するディスプレイ装置において、冗長データ生成器へのデータ入力、及び、冗長データ生成器からのデータ出力の、転送効率の向上を可能とする。

【0035】また、データリード時、冗長データも含めてリードするため、冗長データの信頼性を向上させることが可能となる。

【0036】本発明によれば、データライト時、冗長データ生成器により、上位装置からのデータはデータ分散単位を分割した単位で入力可／不可制御され、また、生成した冗長データは、同様にデータ分散単位を分割した単位で出力制御される。このため、冗長データは、分散単位分の冗長データが生成されるのを待つことなく、1つの分割単位の冗長データが生成される毎に出力し、また、上位装置からのデータは、データ分散単位×分散ディスク装置数のデータ転送を行った後、生成した分散単位の冗長データが全て冗長データ生成器より出力されるのを待つことなく、分割単位での冗長データ出力に応じた次の分散のデータ入力を行うことができる。

【0037】すなわち、冗長データ生成器へのデータの入出力可／不可制御を、データ分散単位を分割した単位で行うことで、生成した冗長データの出力待ちによるデータ入力不可時間を低減し、データ転送の効率を向上させる。

【0038】また、各々のディスク装置ユニット列の途中の列にてデータライトが開始、あるいは、終了するような場合、この分散データに対する冗長データを生成する場合は、上位装置、及びディスク装置からの冗長データ生成器へのデータ入力が必要であるが、この時、上位装置、各ディスク装置からのデータ入力順、及び、タイミングに制限がないため、データ入力にソフト介在の必要がなく、また、上位装置とディスク装置と同時に転送起動可能となり、ディスク装置からのデータ入力のオーバヘッドが低減できる。

【0039】すなわち、冗長データ生成器へのデータの入出力可／不可制御を、上位装置側、及び、それぞれのディスク装置側全て独立に制御することで、冗長データ生成器へのデータ入力タイミング、データ入力順（上位装置側及びディスク装置側のどちらを先に入力させるかという順序）等を意識する必要がなくなり、またそのためのソフト制御も必要とせず、データ転送の効率を向上させる。

【0040】更に、分散単位の途中にてデータライトが開始、あるいは終了するような場合に、分散単位の開始から途中まで、あるいは、分散単位の途中から終了までをダミー転送することで、上位装置、及び、ディスク装置共にデータ転送の単位は分散単位となり、冗長データ生成のための転送を、分散単位の開始から途中までの列、あるいは、分散単位の途中から終了までの列、及び、それ以外の列とに分割して制御する必要がなくなる。

【0041】すなわち、従来ホストマシンから送られてくるデータの単位よりも分散単位が大きいときに分散単位分のデータをディスク装置ユニットから実際に読み込

まなければならなかったが、ダミーデータを使用することにより、実際に読み込まなければならないデータを減らすことができる。従って、冗長データの生成効率が向上する。

【0042】また、データリード時に冗長データも含めてリードし、冗長データ生成器にてデータ照合するため、データ再構築時の信頼性を向上できる。

【0043】次に、第2の実施例について説明する。

【0044】図3は、本発明を説明する為の、ディスクアレイ装置の実施例である。本実施例の装置は、コンピュータの記憶装置として動作する。本実施例では、汎用I/OインターフェースであるホストSCSIバス11を介して、コンピュータと接続される。本装置の機能は、コンピュータの要求に対して、該当するディスク装置ユニットに、データ書き込み及び、読みだしを行なうことであるが、本機能を実現する手段に、データの分散集合手法を用いていることが第1の特徴としてあげられる。

【0045】本ディスクアレイ装置は、ホストSCSIバス11と、インタフェース制御回路12と、データバス13と、データ分散集合制御回路14と、共通データバス15と、データバッファ171と、DMA回路172と、冗長データ生成及びデータ再構築回路161と、ディスク制御回路18と、ディスク装置ユニット（ディスク駆動機構）19とを有する。

【0046】本実施例においては、コンピュータのデータ書き込み要求時、ホストSCSIバス11を介して転送されるデータは、インタフェース制御回路12を経由し、データバス13よりデータ分散集合制御回路14に入力される。

【0047】ここでコンピュータの指示した論理ユニットアドレス、論理ブロックアドレス、転送ブロック数より、どのディスク装置ユニットのどの位置に、どれだけのデータ転送を行なうか等の、分散パラメータが生成され、コンピュータからのデータが分散される。

【0048】分散されたデータは、共通データバス15を介して、分散先のディスク装置ユニット19を制御するディスク制御回路18に開放されたデータバッファ171のエリアにDMA回路172を介して転送される。各々のディスク制御回路は、自身に開放されたデータバッファのエリアより、指定されたディスク装置ユニット19に対して、分散されたデータを書き込む。

【0049】コンピュータの、データ読みだし要求時は、コンピュータの指示するパラメータ及び分散情報より、該当するディスク制御回路181が、指定されたディスク装置ユニット19よりデータを読みだし、データバッファ171の定められたエリアに転送する。データ分散集合制御回路14は、データバッファ171に転送された分散データを集合し、集合されたデータは、データバス13、インタフェース制御回路を経由して、ホス

トSCSIバス11を介してコンピュータに転送される。

【0050】本ディスクアレイ装置の第2の特徴として、データの冗長性があげられる。ディスクアレイ装置に冗長性を持たず場合は、コンピュータのデータ書き込み要求時、データ分散集合制御回路14により分散されるデータを、データバッファ171に転送すると同時に、冗長データ生成及びデータ再構築回路16に転送する。転送された各々の分散データは、分散するディスク装置ユニット数分、排他ORが実行されて、生成回路内に冗長データが生成される。

【0051】本冗長データの転送先は、冗長データ生成及びデータ再構築回路16にデータが入力される際に、データ分散集合制御回路141より、冗長データ生成及びデータ再構築回路16に転送されているため、冗長データの生成が完了すると、分散データと同様にして該当するデータバッファ171のエリアに転送され、ディスク制御回路181を経由してディスク装置ユニットに書き込まれる。

【0052】このように、分散データから生成され、ディスク装置ユニットに書き込まれた冗長データは、分散を行なった複数のディスク装置ユニットの中で、あるディスク装置ユニットに障害が発生した場合、障害の発生したディスク装置ユニットの分散データを復元するため

に使用される。

【0053】本装置では、障害の発生したディスク装置ユニット以外のディスク装置ユニットより、分散されたデータ及び、冗長データを読みだし、データは、該当するデータバッファ171内のエリアと、冗長データ生成及びデータ再構築回路16へ転送し、冗長データは、冗長データ生成及びデータ再構築回路16へのみ転送することで、冗長データ生成及びデータ再構築回路16は、障害分の分散データを再構築する。

【0054】この再構築されたデータは、該当するデータバッファ171のエリアに転送される。データ分散集合制御回路141は、通常の読みだし処理と同様に、データバッファ171の分散データを集合し、集合されたデータは、コンピュータへと転送される。

【0055】また、ディスクアレイ装置内に、スペアディスク装置ユニットを備えている場合は、ディスクアレイ装置内で、障害ディスク装置ユニットの内容を、スペアディスクに再構築する。

【0056】以上、ディスクアレイ装置の実施例を説明したが、本発明は、コンピュータのデータを、データ分散、集合手段を用いて装置内で加工する為のデータ保証方法であり、その実施例を図4、図5に示す。

【0057】図4は、データ分散集合制御回路141のブロック図を表わす。本データ分散集合制御回路141は、分散レングスカウント121と、分散パラメータ発生器（分散パラメータ付加手段）122と、バッファ分

散アドレス生成器123と、分散パラメータ照合器（分散パラメータ確認手段）124と、FIFOメモリ123とを有する。

【0058】データバス13より、入力されたコンピュータからのデータは、分散レングスカウンター121により、最小分散単位の512バイトのデータ長に分割され、FIFOメモリ125に入力される。このとき、分散されたデータに、分散パラメータ発生器122によって生成された分散パラメータ、本実施例では、8バイト長が、最小分散データ単位に付加されてFIFOメモリ125に入力される。FIFOメモリ125に入力された各々の分散データ及び分散パラメータは、共通バス15を介して、データバッファ171及び、冗長データ生成及びデータ再構築回路116へDMA転送される。このときのデータバッファ171転送アドレスは、バッファ分散アドレス生成器123によって生成される。

【0059】コンピュータにデータを転送する際は、データバッファ171に転送される分散データを、バッファ分散アドレス生成器123によって生成されるバッファアドレスよりFIFOメモリ25に転送する。FIFOメモリ25に転送された各々の分散データ及び、分散パラメータは、FIFOメモリ25上に分散された順序で集合しており、FIFOメモリ25よりデータバス13にデータ転送をする際、分散パラメータ照合器24が分散パラメータを照合して、分散データより切り離し、データのみをデータバス13に転送する。

【0060】図5は、ディスク制御回路181のブロッ

表1 分散パラメータ構成

付加パラメータ名	ビット長	内 容
ディスク制御番号	4	ディスク装置ユニットを制御するディスク制御回路の固有番号。
論理ユニット番号	6	上位より指定の論理ユニットアドレス。
シーケンス番号	6	1回で分散する分散単位内のブロックの順位。
パリティグループ番号	32	同一論理ユニット内の（分散単位）×（分散ディスク装置）分のデータを管理する番号。
LRCコード	16	上記パラメータを含めた分散ブロック単位の冗長コード。

【0067】なお、本実施例において、分散パラメータに規則性を持たせる、すなわち規則的にデータを分散させると、分散パラメータの排他論理和をとったものも予

測可能なものとなる。

【0061】本ディスク制御回路181は、FIFOメモリ31と、分散パラメータ照合器32と、SCSI制御器33と、ディスクSCSIバス34とを有する。

【0062】共通データバス15を介して入力される各々の分散データ及び、分散パラメータは、FIFOメモリ31に転送される。FIFOメモリ31に転送された分散データ及び分散パラメータは、分散パラメータ照合器32で分散データの妥当性が照合され、SCSI制御器33を介して、該当のディスク装置ユニットに、分散パラメータも含めて書き込まれる。

【0063】ディスクSCSIバス34より入力される分散データ及び分散パラメータは、SCSI制御器を介してFIFOメモリ31へ転送される。このとき分散パラメータは、分散パラメータ照合器32の妥当性が照合する。FIFOメモリ31に転送された分散データ及び分散パラメータは、共通データバスへ出力される。

【0064】以上の様に、本発明は、ディスクアレイ装置において、データの分散、集合をおこなう際のデータの保証を可能とする。また、分散パラメータがディスク装置ユニットに書き込まれることから、分散パラメータを読み込むことで、分散情報を再生することが可能となる。

【0065】本実施例における分散パラメータの内容を下記の表1に示す。

【0066】

【表1】

測可能なものとなる。これを冗長データの分散パラメータとみなして、冗長データに付加して記録し、再生時にこの分散パラメータが予測するものになっているかどうか

かを判定することにより、冗長データを記録しているディスク装置ユニットに故障が生じたかがわかる。従って、ディスクアレイ装置の信頼性が向上する。

【0068】本実施例によれば、データ記録時、最小分散長である512バイトのデータ単位に分散パラメータが生成、付加され、分散先のディスク装置ユニット毎に、付加された分散パラメータが照合され、分散データの妥当性を保証して分散データと共にディスク装置に記録される。また冗長データにおいても、分散パラメータ分の冗長データが生成され、冗長データとしての妥当性を照合し、ディスク装置に記録される。データ読みだし時は、ディスク装置に記録されたデータと共に分散パラメータが読みだされ、集合時にデータの妥当性を保証して、正当なデータが上位システムに転送される。また、障害データを他のデータと冗長データより再構築する場合においても、障害データの分散パラメータが他のデータの分散パラメータと冗長データより再構築されるため、集合時にデータの妥当性を保証することが可能となる。

【0069】さらに、分散パラメータを、ディスク装置上に記録するため、ディスクアレイ装置内に保持している分散情報が消失した場合においても、ディスク装置が正常であれば、ディスク装置に記録された分散パラメータより分散情報の再構築が可能となる。

【0070】さらに、以下のような効果もある。

【0071】通常、コンピュータから送られるデータは、ブロック単位（分散単位よりも小さい単位。例えば、ホストシステムからの転送単位が512バイトで、分散単位が2Kバイトのときの512バイト）にチェックコード等がディスクアレイ装置側により付加され、読みだし時に照合されるが、ディスクアレイ装置では、コンピュータからのデータを装置内で分散集合加工するため、ブロック単位のチェックコード等は正常でも、ブロック単位の分散単位内の順序（シーケンス）が故障により保証されなくなる。分散パラメータの中にシーケンス情報を含ませることにより、シーケンスもチェックすることができる。

【0072】また、ディスクアレイ装置内で生成される

冗長データにも分散パラメータを付加し、冗長データより、障害データを再構築する際、再構築したデータの分散パラメータをも再構築し、分散パラメータより故障したディスク装置ユニットに分散されたものであるかがわかる。もし、故障したディスク装置ユニット以外に分散されたものであるという分散パラメータが検出されたときは、ユニットの故障以外が考えられることにより、再構築されたデータの信頼性が向上する。

【0073】さらに、ディスクアレイ装置内に保持する分散情報が破壊されても、各々のディスク装置ユニットの分散パラメータが正常であれば、ディスク装置ユニットに記録された分散パラメータより、分散情報を再生して、データ破壊を防止する。

【0074】

【発明の効果】本発明によれば、冗長データのディスク装置への転送時間を短縮したディスクアレイ装置を提供できる。

【0075】また、障害検出の遅れを防ぎ、信頼性を向上させたディスクアレイ装置を提供できる。

【0076】また、分散、集合に対する信頼性を向上させたディスクアレイ装置を提供できる。

【図面の簡単な説明】

【図1】冗長データを有するディスクアレイ装置のブロック図である。

【図2】冗長データ生成器のブロック図である。

【図3】ディスクアレイ装置のブロック図である。

【図4】データ分散集合制御回路のブロック図である。

【図5】ディスク制御回路のブロック図である。

【符号の説明】

11…ホストSCSIバス、12…インターフェース制御回路、13…データバス、14…データ分散集合制御回路、15…共通データバス、16…冗長データ生成器、17…データバッファ、18…ディスク制御回路、19…ディスク装置ユニット、21…上位側データ管理カウンタ、22…ディスク側データ管理カウンタ、23…冗長データ管理カウンタ、24…入出力制御回路、25…入出力可／不可信号。

ディスクアレイ装置ブロック図(図1)

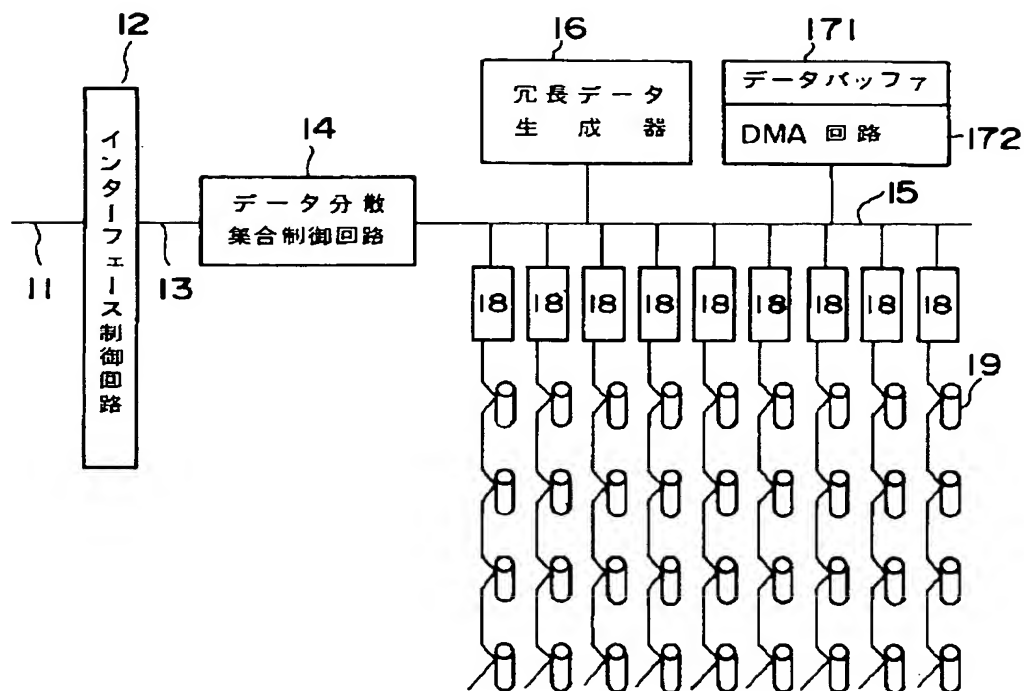
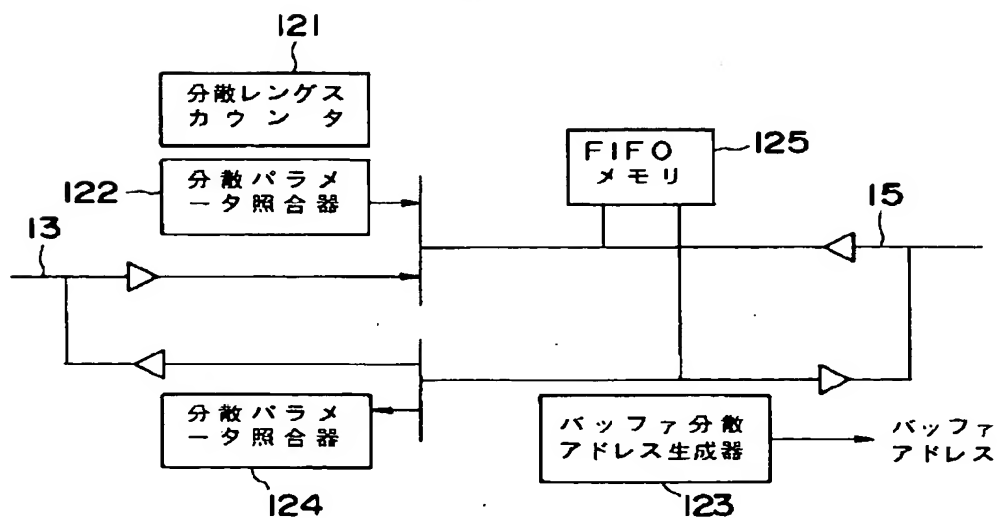
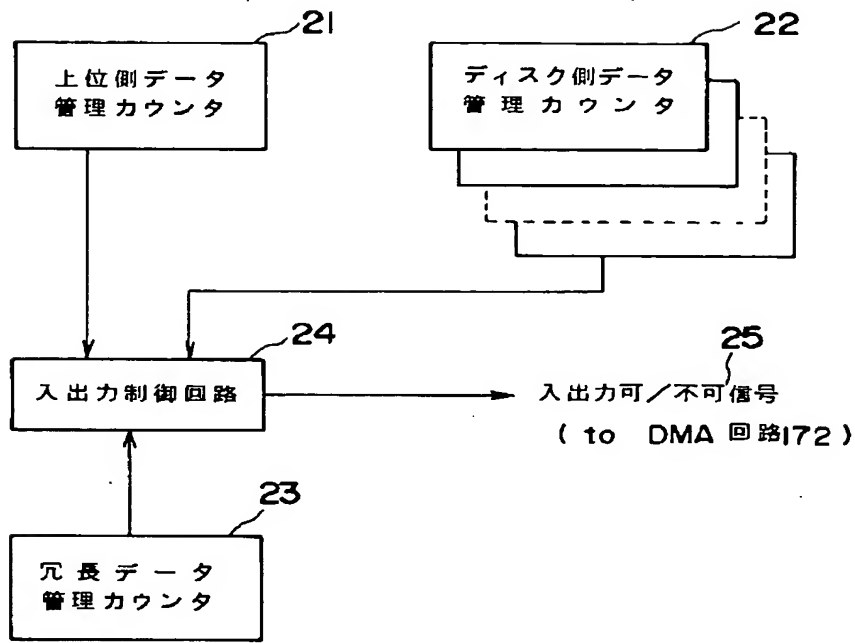


Figure 4 shows a schematic diagram of a rectangular domain with a central square hole. The domain is divided into four quadrants by a vertical line and a horizontal line. The central square hole is also divided into four quadrants. The outer boundary is labeled 'a' and the inner boundary is labeled 'b'. The domain is labeled 'D'.



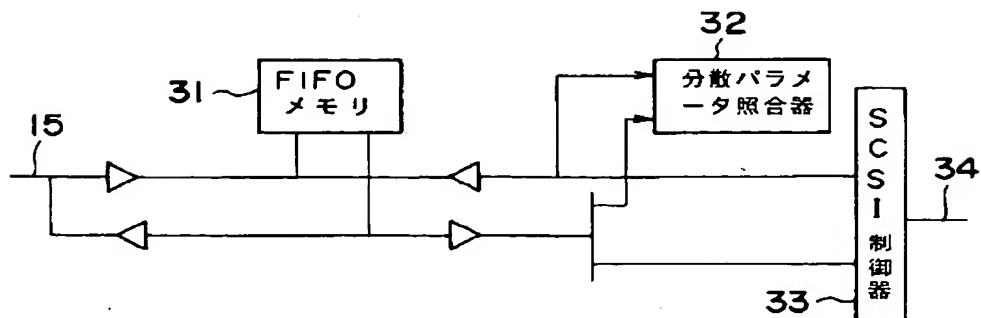
【図2】

冗長データ生成器入出力制御ブロック図（図2）



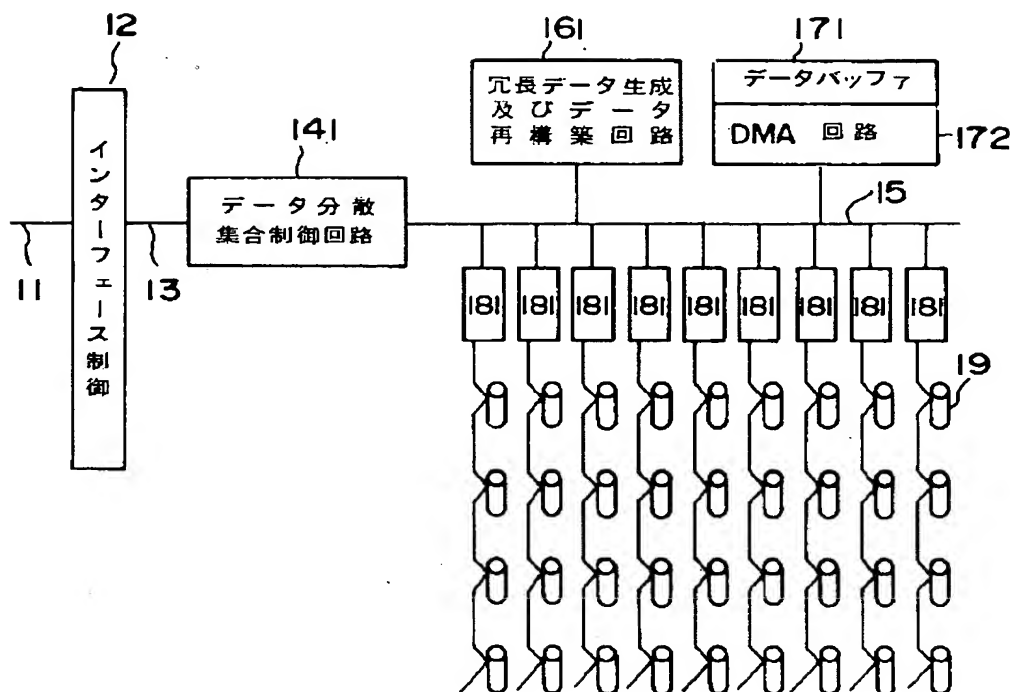
【図5】

図 5



【図3】

ディスクアレイ装置ブロック図（図3）



フロントページの続き

(72)発明者 水野 克敏
神奈川県小田原市国府津2880番地 日立コ
ンピュータ機器株式会社内

35

(72)発明者 鈴木 良一
神奈川県小田原市国府津2880番地 日立コ
ンピュータ機器株式会社内
(72)発明者 馬場 英美
神奈川県小田原市国府津2880番地 日立コ
ンピュータ機器株式会社内